

INTRODUCTION TO DEEP LEARNING

Winter School at UPC TelecomBCN Barcelona. 22-30 January 2018.



Instructors



Xavier
Giró-i-Nieto

Marta R.
Costa-jussà

Noé
Casas

Elisa
Sayrol

Antonio
Bonafonte

Verónica
Vilaplana

Ramon
Morros

Javier
Ruiz

Organizers



UNIVERSITAT POLITÈCNICA
DE CATALUNYA
BARCELONATECH



Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

Supporters

aws  educate

GitHub Education

+ info: <https://telecombcn-dl.github.io/2018-idl/>

[\[course site\]](#)



#DLUPC

Day 2 Lecture 3

Loss functions



Javier Ruiz Hidalgo

javier.ruiz@upc.edu

Associate Professor

Universitat Politècnica de Catalunya
Technical University of Catalonia



Me



- **Javier Ruiz Hidalgo**

- Email: j.ruiz@upc.edu
- Office: UPC, Campus Nord, D5-008

- **Teaching experience**

- Basic signal processing
- Project based
- Image processing & computer vision

- **Research experience**

- Master on hierarchical image representations by UEA (UK)
- PhD on video coding by UPC (Spain)
- Interests in image & video coding, 3D analysis and super-resolution

Outline

- **Introduction**
 - **Definition, properties, training process**
- **Common types of loss functions**
 - Regression
 - Classification

Definition

In a supervised deep learning context the **loss function** measures the **quality** of a particular set of parameters based on how well the output of the network **agrees** with the ground truth labels in the training data.

Nomenclature

loss function

=

cost function

=

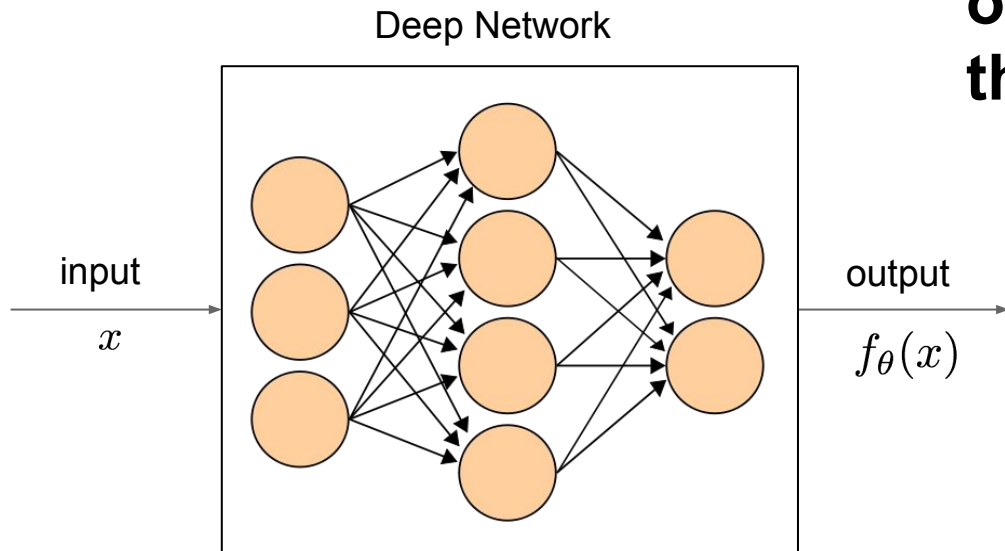
objective function

=

error function

Loss function (1)

**How good does
our network with
the training data?**



$$\mathcal{L}(w) = \text{distance}(f_{\theta}(x), y)$$

input → x

labels (ground truth) → y

error → $\mathcal{L}(w)$

parameters (weights, biases) → θ

Loss function (2)

- The loss function **does not** want to **measure** the **entire performance** of the network against a validation/test dataset.
- The loss function is used to **guide** the **training process** in order to find a set of parameters that reduce the value of the loss function.

Training process

Stochastic gradient descent

- Find a set of parameters which make the loss as small as possible.
- Change parameters at a rate determined by the partial derivatives of the loss function:

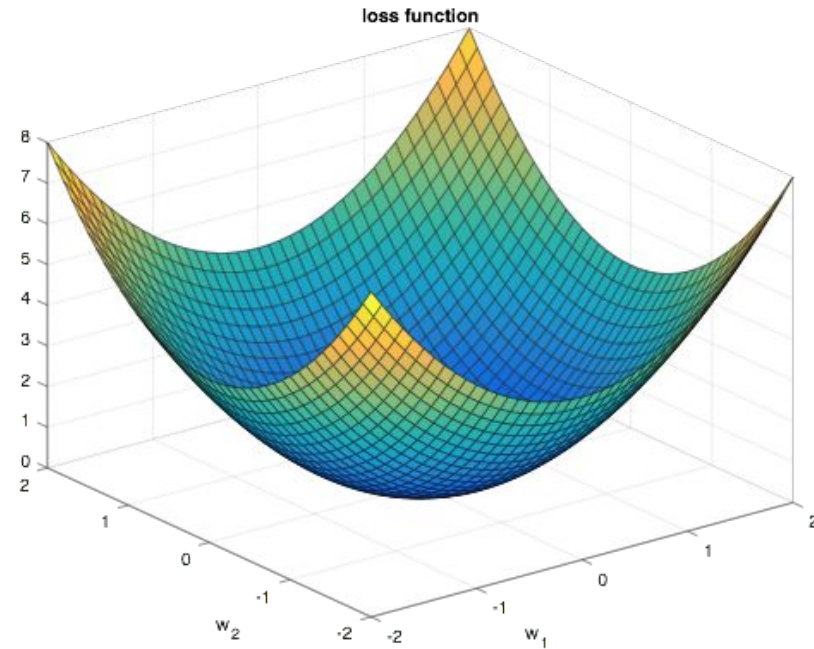
$$\frac{\partial \mathcal{L}}{\partial w} \quad \frac{\partial \mathcal{L}}{\partial b}$$

Properties (1)

- Minimum (0 value) when the output of the network is equal to the ground truth data.
- Increase value when output differs from ground truth.

Properties (2)

- Ideally \rightarrow convex function
- In reality \rightarrow many parameters (in the order of millions) not convex
- Varies smoothly with changes on the output
 - Better gradients for gradient descent
 - Easy to compute small changes in the parameters to get an improvement in the loss



Outline

- Introduction
 - Definition, properties, training process
- **Common types of loss functions**
 - **Regression**
 - Classification

Common types of loss functions (1)

- Loss functions depend on the type of task:
 - Regression: the network predicts **continuous, numeric** variables
 - Example: Length of fishes in images, temperature from latitude/longitude
 - Absolute value, square error

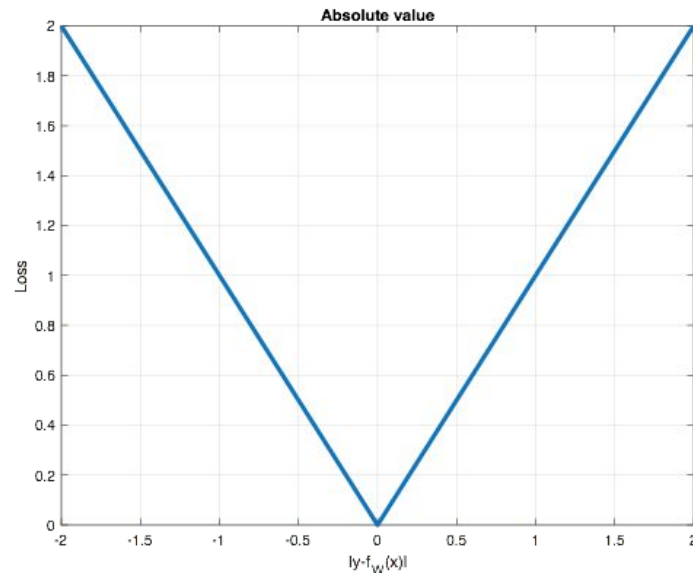
Common types of loss functions (2)

- Loss functions depend on the type of task:
 - Classification: the network predicts **categorical** variables (fixed number of classes)
 - Example: classify email as spam, predict student grades from essays.
 - hinge loss, Cross-entropy loss

Absolute value, L1-norm

- Very intuitive loss function
 - produces sparser solutions
 - good in high dimensional spaces
 - prediction speed
 - less sensitive to outliers

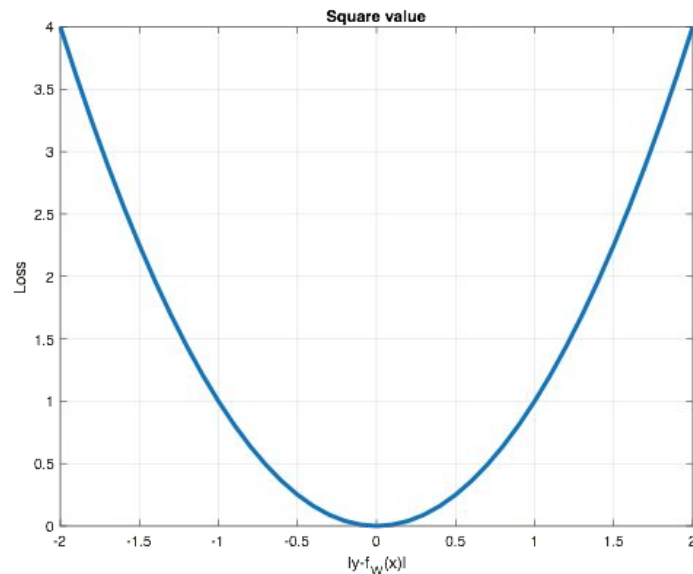
$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n |y_i - f_{\theta}(x_i)|$$



Square error, Euclidean loss, L2-norm

- Very common loss function
 - More precise and better than L1-norm
 - Penalizes large errors more strongly
 - Sensitive to outliers

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n (y_i - f_{\theta}(x_i))^2$$

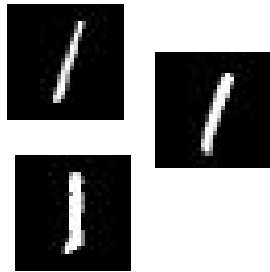


Outline

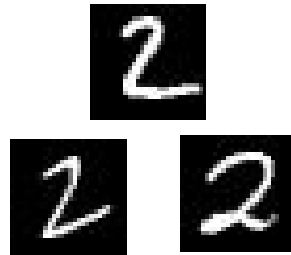
- Introduction
 - Definition, properties, training process
- **Common types of loss functions**
 - Regression
 - **Classification**

Classification (1)

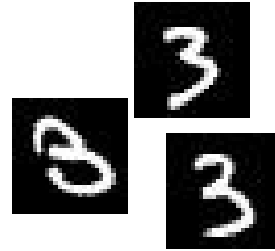
We want the network to classify the input into a fixed number of classes



class "1"



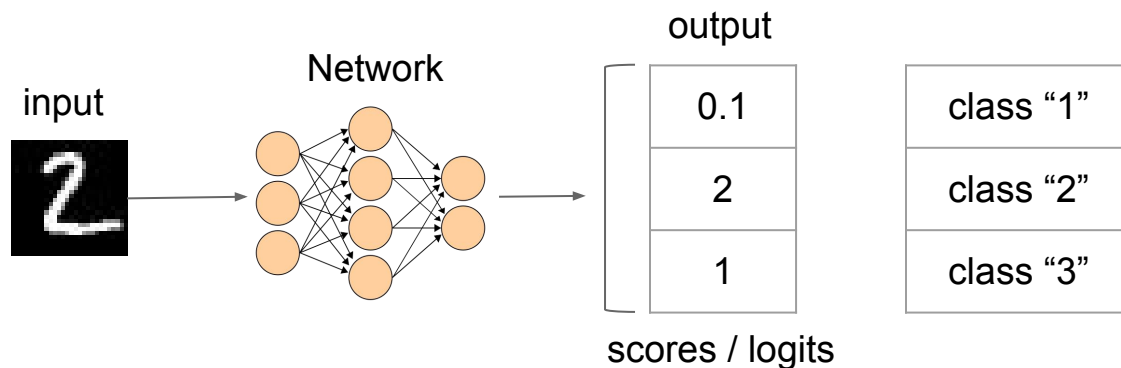
class "2"



class "3"

Classification (2)

- Each input can have only one label
 - One prediction per output class
 - The network will have “k” outputs (number of classes)



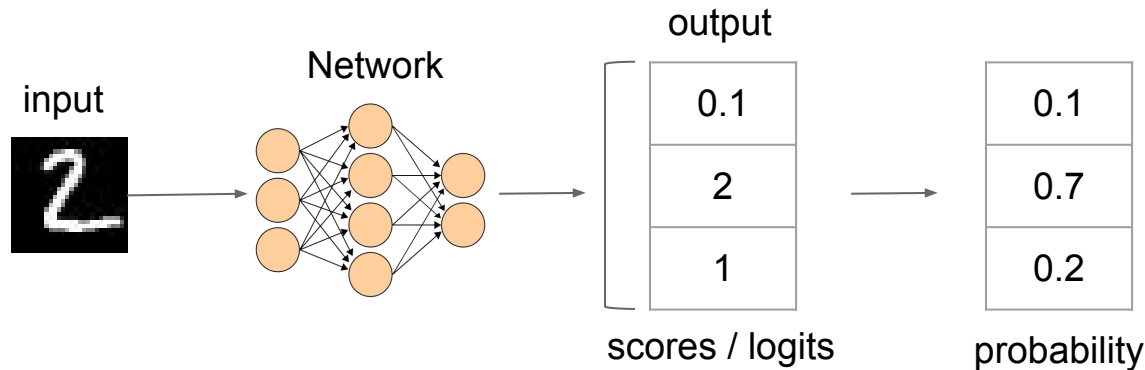
Classification (3)



- How can we create a loss function to improve the scores?
 - Somehow write the labels (ground truth of the data) into a vector → One-hot encoding
 - Non-probabilistic interpretation → **hinge loss**
 - Probabilistic interpretation: need to transform the scores into a probability function → Softmax

Softmax (1)

- Convert scores into probabilities
 - From 0.0 to 1.0
 - Probability for all classes adds to 1.0

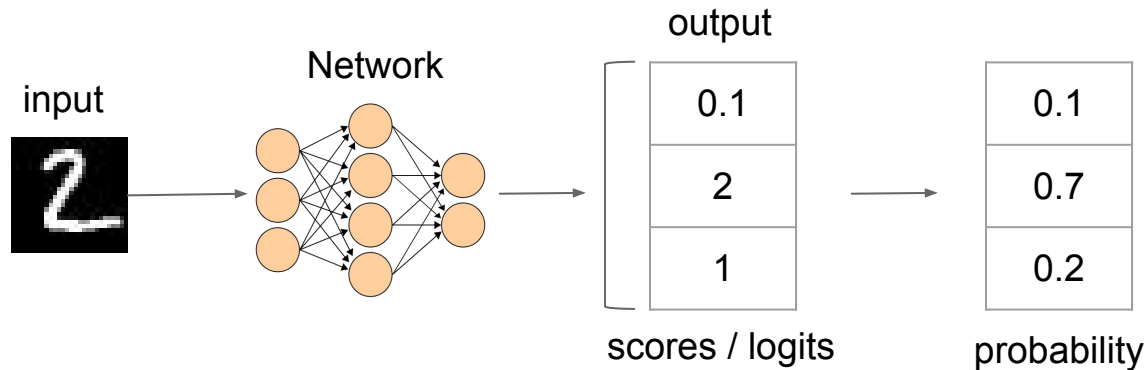


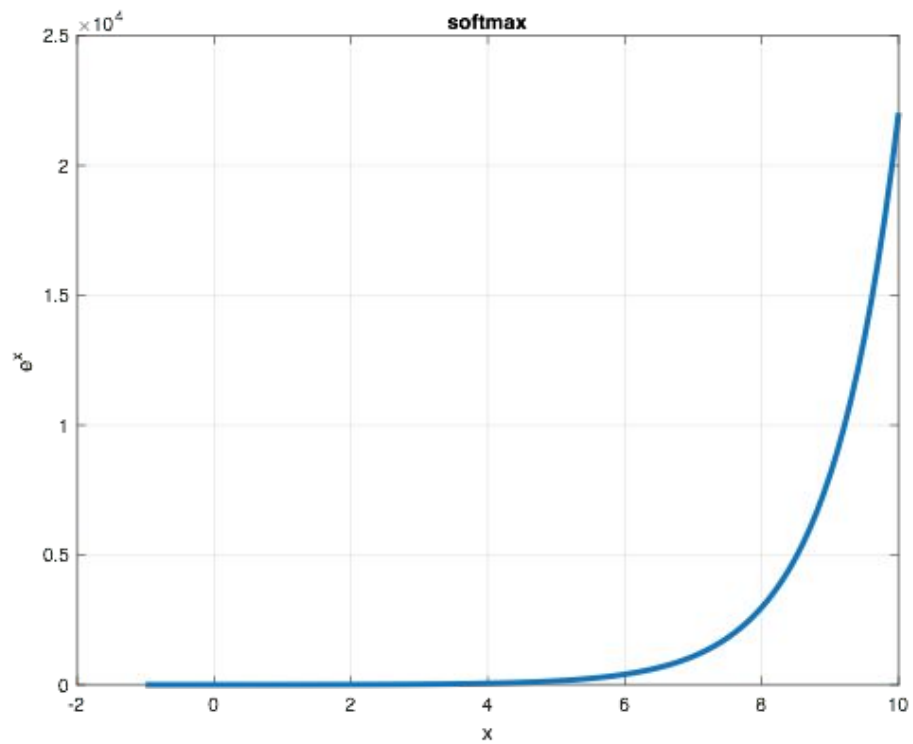
Softmax (2)

- Softmax function

scores (logits)

$$S(l_i) = \frac{e^{l_i}}{\sum_k e^{l_k}}$$



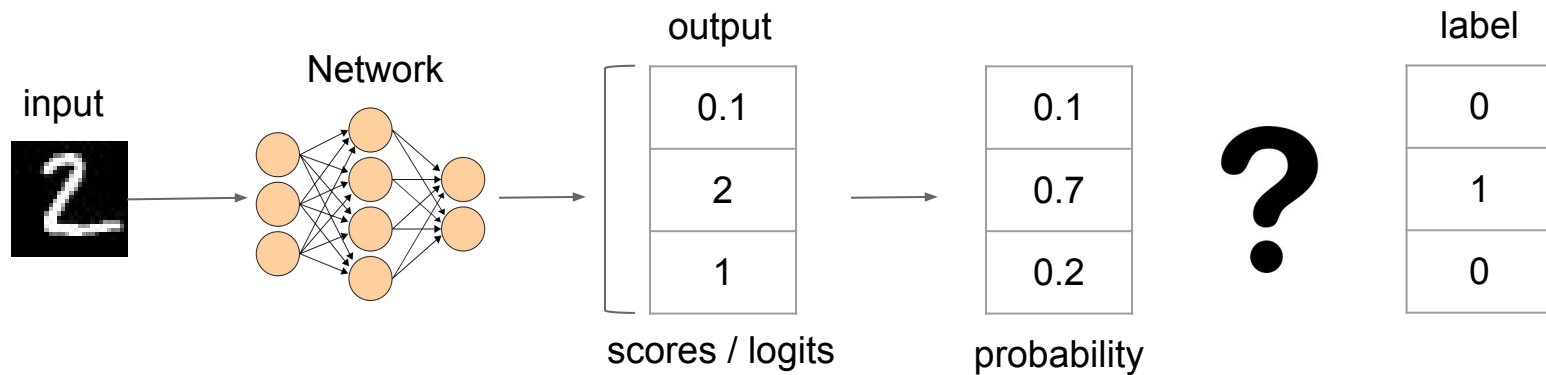


One-hot encoding

- Transform each label into a vector (with only 1 and 0)
 - Length equal to the total number of classes “k”
 - Value of 1 for the correct class and 0 elsewhere

class “1”	class “2”	class “3”
1	0	0
0	1	0
0	0	1

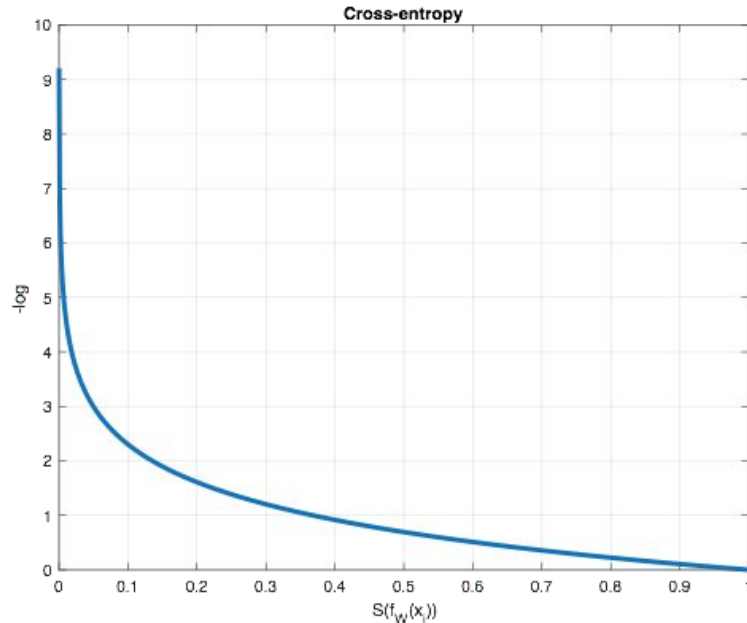
Cross-entropy loss (1)



$$\mathcal{L}_i = - \sum_k y_k \log(S(l_k)) = - \log(S(l))$$

Cross-entropy loss (2)


$$\mathcal{L}_i = - \sum_k y_k \log(S(l_k)) = -\log(S(l))$$




Cross-entropy loss (3)

- For a set of n inputs $\mathcal{L} = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i$

$$\mathcal{L} = - \sum_{i=1}^n \mathbf{y}_i \log(S(f_{\theta}(\mathbf{x}_i)))$$

labels (one-hot) 

Softmax 


Cross-entropy loss (4)

- In general, cross-entropy loss works better than square error loss:
 - Square error loss usually gives too much emphasis to incorrect outputs.
 - In square error loss, as the output gets closer to either 0.0 or 1.0 the gradients get smaller, the change in weights gets smaller and training is slower.

Regularization

- Control the capacity of the network to **prevent overfitting**

- L2-regularization (weight decay): regularization parameter

$$\mathcal{L}_{new} = \mathcal{L} + \frac{\lambda}{2} W^2$$


- L1-regularization:

$$\mathcal{L}_{new} = \mathcal{L} + \frac{\lambda}{2} |W|$$

References

- [About loss functions](#)
- [Neural networks and deep learning](#)
- [Are loss functions all the same?](#)
- [Convolutional neural networks for Visual Recognition](#)
- [Deep learning book, MIT Press, 2016](#)
- [On Loss Functions for Deep Neural Networks in Classification](#)

Thanks! Questions?



UNIVERSITAT POLITÈCNICA DE CATALUNYA
BARCELONATECH

Department of Signal Theory
and Communications

Image Processing Group

<https://imatge.upc.edu/web/people/javier-ruiz-hidalgo>